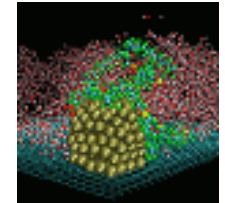
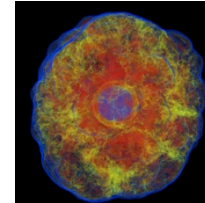
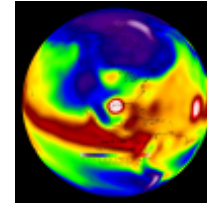
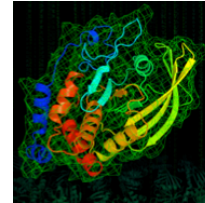
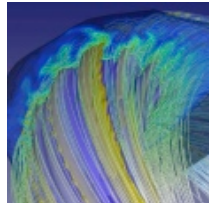
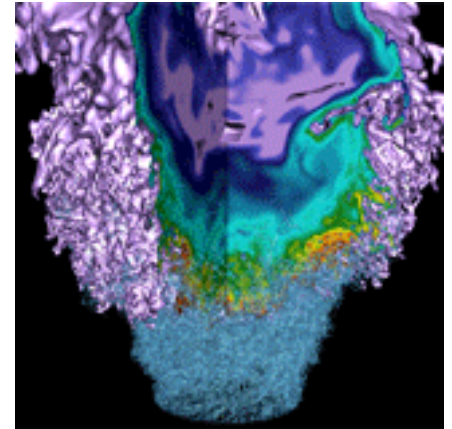


# Data and Analytics Strategy



**Prabhat**

Data and Analytics Group Lead  
February 23, 2015

# Talk Overview

---

- **DAS Team and Goals**
- **Big Data Hardware**
- **Big Data Software**
- **Big Data Users**

# Talk Overview

---

- **DAS Team and Goals**
- Big Data Hardware
- Big Data Software
- Big Data Users

# Data and Analytics Team

DAS Team Member	Technology Areas
Shreyas Cholia	Gateways, Web, Grid
Yushu Yao	Databases, Analytics
Annette Greiner	UI, Web,
Joaquin Correa	Imaging, Machine Learning
Burlen Loring	Vis
Jeff Porter	Data Management
Oliver Ruebel	Vis, Analytics
Dani Ushizima	Imaging, R
R. K. Owen	NIM
Michael Urashka	Web

# DAS Team Goal: “Enable Data-Centric Science at Scale”

---

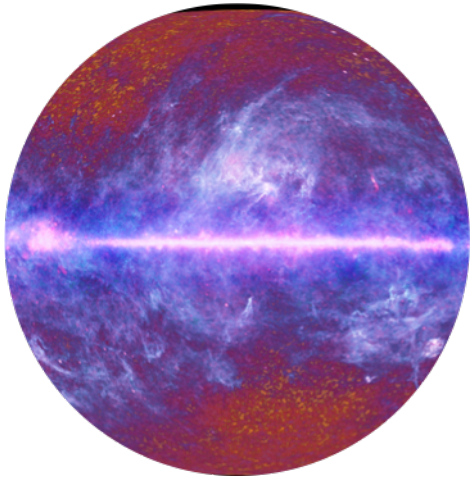
- **Big Data Software**
  - Broad ecosystem of capabilities and technologies
  - Research and evaluate
  - Customize and optimize for NERSC/HPC platforms
  - Deploy and maintain
- **Engaging NERSC Users**
  - Broad user base support
  - 1-1 in-depth engagement

# Talk Overview

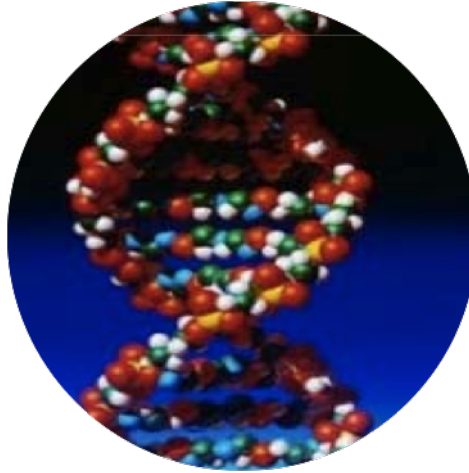
---

- DAS Team and Goals
- **Big Data Hardware**
- Big Data Software
- Big Data Users

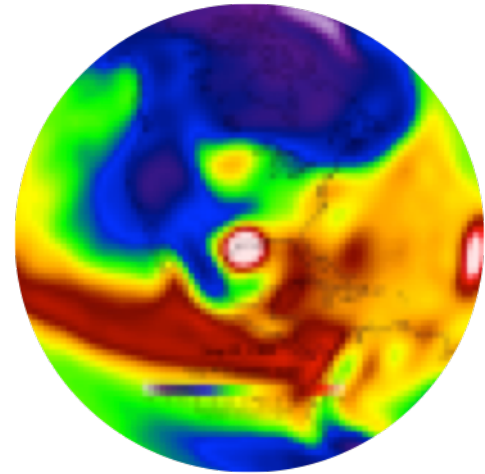
# DOE Facilities are Facing a Data Deluge



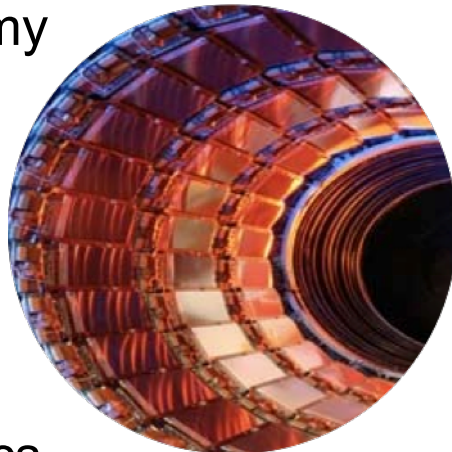
Astronomy



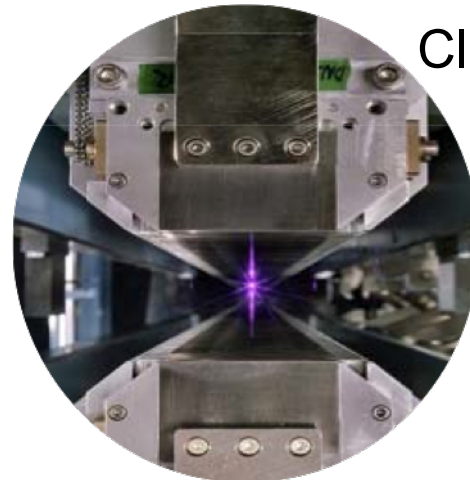
Genomics



Climate



Physics



Light Sources



# We currently deploy separate Compute Intensive and Data Intensive Systems

## *Compute Intensive*



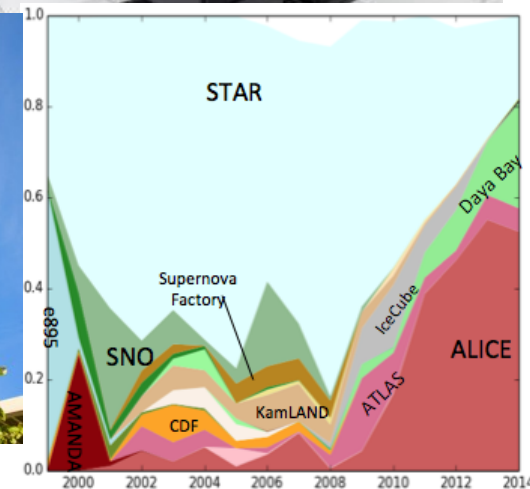
## *Data Intensive*



Carver



Genepool



PDSF



# Cori: Unified architecture for HPC and Big Data

---

- **64 Cabinets of Cray XC System**
  - 50 cabinets ‘Knights Landing’ *manycore* compute nodes
  - 10 cabinets ‘Haswell’ compute nodes for *data partition*
  - ~4 cabinets of Burst Buffer
  - 14 external login nodes
  - Aries Interconnect (same as on Edison)
- **Lustre File system**
  - 28 PB capacity, 432 GB/sec peak performance
- **NVRAM “Burst Buffer” for I/O acceleration**
- **Significant Intel and Cray application transition support**
- **Delivery in mid-2016; installation in new LBNL CRT**

# Popular features of a data intensive system can be supported on Cori

Data Intensive Workload Need	Cori Solution
Local Disk	NVRAM 'burst buffer'
Large memory nodes	128 GB/node on Haswell; Option to purchase fat (1TB) login node
Massive serial jobs	NERSC serial queue prototype on Edison; MAMU
Complex workflows	More (14) external login nodes; CCM mode for now
Communicate with databases from compute nodes	<b><i>Proposed Compute Gateway Node COE</i></b>
Stream Data from observational facilities	<b><i>Proposed Compute Gateway Node COE</i></b>
Easy to customize environment	<b><i>Proposed User Defined Images COE</i></b>
Policy Flexibility	Improvements coming with Cori: Rolling upgrades, CCM, MAMU, above COEs would also contribute

# Talk Overview

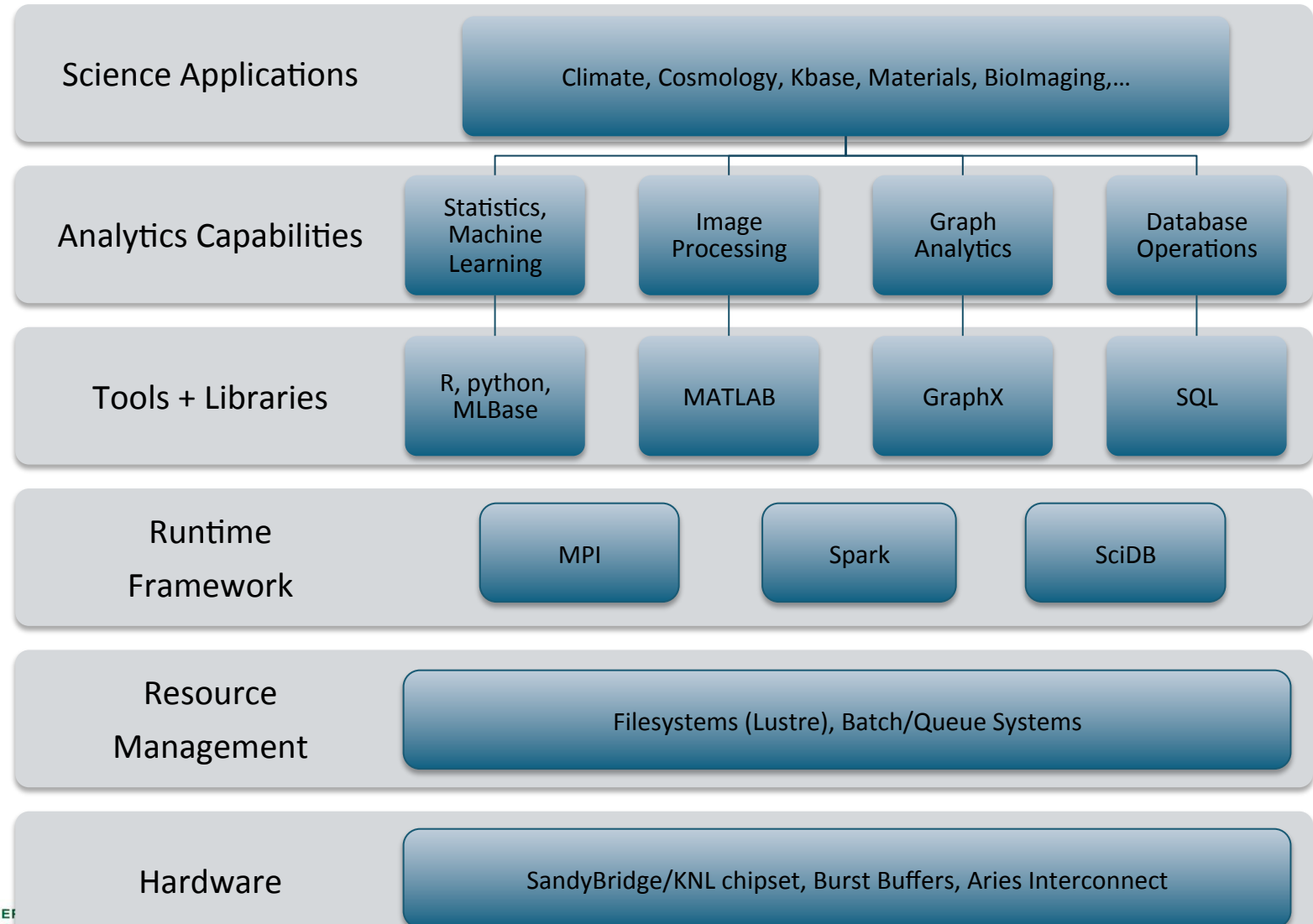
---

- DAS Team and Goals
- Big Data Hardware
- **Big Data Software**
- Big Data Users

# Big Data Software Portfolio

Capabilities	Technology Areas	Tools, Libraries
Data Transfer + Access	Globus, Grid Stack, Authentication	Globus Online, Grid FTP
	Portals, Gateways, RESTful APIs	NEWT
Data Processing	Workflows	Swift, Fireworks, ...
Data Management	Formats, Models Databases	HDF5, NetCDF
	Storage, I/O, Movement	SRM
Data Analytics	Statistics, Machine Learning	python, R, ROOT
	Imaging	OMERO, Fiji, ...
Data Visualization	SciVis InfoVis	VisIt, Paraview
Backend Infrastructure	Analytics Stack Databases  Virtualization	BDAS SciDB, MySQL, PostgreSQL, MongoDB Docker

# Analytics Software Strategy



# Current DAS Engagements

---

- **Analytics:**
  - Cray, UCB AMPLab, Databricks, SkyTree, Dato
  - Intel Research, Nervana Systems, UCB, Harvard, MIT, CMU
- **Data Transfer, Access:**
  - Globus
- **Visualization**
  - Kitware
- **Data Management:**
  - HDF Group
  - Paradigm4, MongoDB

# Talk Overview

---

- DAS Team and Goals
- Big Data Hardware
- Big Data Software
- **Big Data Users**



# NERSC Users: How to get help?

---

- **Documentation:**
  - <http://www.nersc.gov/users/software/data-visualization-and-analytics/>
- **Routine startup/troubleshooting questions:**
  - Trouble ticket system
- **In-depth 1-1 collaborations:**
  - e-mail [prabhat@lbl.gov](mailto:prabhat@lbl.gov)

# Top NERSC Production Workflows

---

- **Advanced Light Source SPOT suite**
  - Real time reconstruction, experimental steering
- **Materials Project**
- **Cosmology Supernovae/Transient classification pipeline**

---

# Questions?

Contact: [prabhat@lbl.gov](mailto:prabhat@lbl.gov)